# General Disclaimer

## One or more of the Following Statements may affect this Document

- This document has been reproduced from the best copy furnished by the organizational source. It is being released in the interest of making available as much information as possible.

- This document may contain data, which exceeds the sheet parameters. It was furnished in this condition by the organizational source and is the best copy available.

- This document may contain tone-on-tone or color graphs, charts and/or pictures, which have been reproduced in black and white.

- This document is paginated as submitted by the original source.

- Portions of this document are not fully legible due to the historical nature of some of the material. However, it is the best reproduction available from the original submission.

DEPARTMENT OF MATHEMATICS

UNIVERSITY OF HOUSTON          HOUSTON, TEXAS

MULTIDIMEN. STOCHASTIC APPROX.
USING LOCALLY CONTRACTIVE FNS.
W. M. LAWTON, INST FOR ADV STUDY
PRINCETON, NEW JERSEY
REPORT #46   AUGUST 1975

3801 CULLEN BLVD.
HOUSTON, TEXAS  77004

Multidimensional Stochastic Approximation Using
Locally Contractive Functions

August, 1975

by

Wayne M. Lawton
Institute for Advanced Study
Princeton, New Jersey

# Bibliography

1. C.C. Blaydon, K.S. Fu, and R.L. Kashyap, "Stochastic approximation", Adaptive, Learning and Pattern Recognition Systems, Academic Press, New York and London, 1970, Edited by J.M. Mendel and K.S. Fu.

2. A.S. Householder, Theory of Matrices and Numerical Analysis, Blaisdell Publishing Co., New York, 1964.

3. B. Charles Peters, Jr. and Homer F. Walker, "An iterative procedure for obtaining maximum-likelihood estimates of the parameters for a mixture or normal distributions", to appear.

# Multidimensional Stochastic Approximation
# Using Locally Contractive Functions

1. <u>Summary</u>. A Robbins-Monro type multidimensional stochastic approximation algorithm which converges in mean square and with probability one to the fixed point of a locally contractive regression function is developed. The algorithm is applied to obtain maximum likelihood estimates of the parameters for a mixture of multivariate normal distributions.

2. <u>Introduction</u>. Let $E_k$ be real $k$-dimensional Euclidean space with inner product denoted by $< , >$ and norm denoted by $\| \cdot \|$. Corresponding to every positive definite real $k \times k$ matrix $B$ we define the <u>B-inner product</u> $<x,y>_B = <x,By>$ and the <u>B-norm</u> $\| \times \|_B = <x,Bx>^{1/2}$, for $x,y \in E_k$. A function $\overline{F}:D \longrightarrow E_k$, where $D$ is an open subset of $E_k$, is <u>locally contractive</u> at a point $\theta^\circ \in D$ if there exists a $B$-norm on $E_k$ and a number $\lambda$, $0 \leq \lambda < 1$ such that

$$(2.1) \qquad \| \theta^\circ - \overline{F}(\theta) \|_B \leq \lambda \| \theta - \theta \|_B$$

whenever $\theta$ is sufficiently near $\theta^\circ$. If the above inequality holds for every $\theta$ in some neighborhood $W$ of $\theta^\circ$, we say $\overline{F}$ is <u>$\lambda$-locally contractive at</u> $\theta^\circ$ <u>throughout W</u>. Clearly, $\theta^\circ$ will be a unique fixed point in $W$ for $\overline{F}$. For any $k \times k$ matrix $A$, let the spectral radius of $A$ be denoted by $\rho(A) = \sup\{|\lambda|:\lambda$ is an eigenvalue of $A\}$. The Frechet derivative of $\overline{F}$, if it exists, will be denoted by $\overline{\nabla F}$. The following result, a consequence of Taylor's theorem and the theory in [2; section 2.3], will be used in part

4 of this paper.

(2.2) <u>Lemma</u> If $\overline{\nabla F}$ exists and is continuous in a neighborhood of $\theta^o$, a necessary and sufficient condition for $\overline{F}$ to be locally contractive at $\theta^o$ is that $\overline{F}(\theta^o) = \theta^o$ and $\rho(\overline{\nabla F}(\theta^o)) < 1$.

Let $\{\overline{Y}(\theta): \theta \in D\}$ be a family of random variables with values in $E_k$ satisfying the following conditions

(2.3)     $\sup_{\theta \in D} E(\overline{Y}(\theta)^2) < \infty$     (E denotes conditional expectation with $\theta$ fixed)

(2.4)     the <u>regression function</u> of $\{\overline{Y}(\theta)\}$, denoted by $\overline{M}(\theta) = E(\overline{Y}(\theta))$, can be expressed as $\overline{M}(\theta) = \theta - \overline{F}(\theta)$ where $\overline{F}(\theta)$ is locally contractive at $\theta^o \in D$.

In part 3 of this paper we develop an algorithm which, given the conditions above and given a sufficiently close approximation to $\theta^o$, yields a sequence of recursively defined random variables with values in $E_k$ which converges to $\theta^o$ in mean square and with probability one.

In part 4 of this paper, this algorithm is used to formulate a stochastic analog of the iterative procedure developed by Peters and Walker in [3] for obtaining maximum likelihood estimates of the parameters for a mixture of normal distributions.

(3)   <u>Derivation of the Algorithm</u>

(3.1) <u>Lemma</u> Let $\rho_1 = (\rho(B^{-1}))^{-1}$ and $\rho_2 = \rho(B)$ where B is a positive definite $k \times k$ matrix, and let $r_1, r_2$ be positive real numbers such that

$\dfrac{r_1}{r_2} \le \dfrac{\rho_1}{\rho_2}$. Let $\theta^o$, $\theta \in E_k$ such that $\|\theta^o - \overline{\theta}\| \le r_1$ and let

$S = \{\theta \in E_k : \|\theta - \overline{\theta}\| \le r_2\}$. Define a function $\varphi : E_k \rightarrow S$ as follows

$$\varphi(\theta) = \begin{cases} \theta & \text{if } \theta \in S \\ (1-t)\overline{\theta} + t\theta & \text{where } t = \dfrac{r_2}{\|\overline{\theta}-\theta\|} \quad \text{if } \theta \notin S \end{cases}$$

Then the following inequality holds for every $\theta \in E_k$

$$\|\theta^o - \varphi(\theta)\|_B^2 \le \|\theta^o - \theta\|_B$$

Proof. We may assume $t = \dfrac{r_2}{\|\overline{\theta}-\theta\|} < 1$ for otherwise $\theta \in S$ implying $\varphi(\theta) = \theta$. By the fundamental theorem of calculus:

$$\|\theta^o - \theta\|_B^2 = \|\theta^o - \varphi(\theta)\|_B^2 + \int_{s=t}^{s=1} \frac{\partial}{\partial s}[\|\theta^o - (1-s)\overline{\theta} - s\theta\|_B^2]ds$$

It suffices to prove that the integrand is nonnegative for all $s$, $t \le S \le 1$. Consider

$$\frac{\partial}{\partial s}[\|\theta - (1-s)\overline{\theta} - s\theta\|_B^2] = 2\langle\overline{\theta} - \theta, \theta^o - (1-s)\overline{\theta} - s\theta\rangle_B$$

$$= 2s\langle\overline{\theta} - \theta, \overline{\theta} - \theta\rangle_B + 2\langle\overline{\theta} - \theta, \theta^o - \overline{\theta}\rangle_B.$$

Since $t \le s \le 1$, by the principle axis theorem for real symmetric matrices,

the left side of the above expression is bounded below by $2 r_2 \|\bar{\theta}-\theta\| \rho_1$.

Similarly, the right hand side of the above expression is bounded, in absolute

value, by $2 r_1 \|\theta-\theta\| \rho_2$. The result follows from the hypothesis concerning

$r_1, r_2, \rho_1$ and $\rho_2$.

For the remainder of part 3 we adopt the notation and assumptions stated

previously. Furthermore, we assume that $\bar{F}$ is locally contractive at $\theta^{\bullet}$

throughout $S$ and that $S \subseteq D$. Define a family $\{Y(\theta): \theta \in E_k\}$ of random

variables by

(3.2)
$$Y(\theta) = \bar{Y}(\mathcal{G}(\theta)) - \mathcal{G}(\theta) + \theta$$

Then the following inequality is valid for every $\varepsilon > 0$.

(3.3) <u>Corollary</u> $\displaystyle \inf_{\varepsilon \le \|\theta-\theta^{\circ}\|_B} E(<\theta-\theta^{\theta}, Y(\theta)>_B) > 0$

<u>Proof.</u> Since $E(Y(\theta)) = \mathcal{G}(\theta) - \bar{F}(\mathcal{G}(\theta)) - \mathcal{G}(\theta) + \theta$, the expression above $=$

$\displaystyle \inf_{\varepsilon \le \|\theta-\theta^{\circ}\|_B} \{<\theta-\theta^{\circ}, \theta-\theta^{\circ}>_B + <\theta-\theta^{\circ}, \theta^{\circ} - \bar{F}(\mathcal{G}(\theta))>_B\}$. The first term above $= \|\theta-\theta^{\circ}\|_B^2$,

and the second term is bounded above in absolute value by

$$\|\theta-\theta^{\bullet}\|_B \|\bar{F}(\mathcal{G}(\theta)) - \theta^{\circ}\|_B$$

$$\le \lambda\|\theta-\theta^{\bullet}\|_B \|\mathcal{G}(\theta) - \theta^{\circ}\|_B \qquad \text{(by (2.1))}$$

$$\le \lambda\|\theta-\theta^{\circ}\|_B^2 \qquad\qquad \text{(by (3.1))}$$

Therefore $\displaystyle \inf_{\varepsilon \le \|\theta-\theta^{\circ}\|_B} E(<\theta-\theta^{\theta}, Y(\theta)>_B) \ge (1 - \lambda)\varepsilon^2 > 0.$

(3.4) <u>Definition</u>. A <u>gain sequence</u> is a sequence $\{a_\ell\}$ of positive numbers satisfying $\sum_{\ell=1}^{\infty} a_\ell = \infty$ and $\sum_{\ell=1}^{\infty} a_\ell^2 < \infty$.

(3.5) <u>Remark</u>. For any $C > 0$, $\{a_\ell = \frac{c}{\ell}\}$ is a gain sequence since

$$\underset{k \to \infty}{\text{limit}} \sum_{\ell=1}^{K} \frac{c}{\ell} = \underset{k \to \infty}{\text{limit}} \; c \log K = \infty \quad \text{and} \quad \sum_{\ell=1}^{\infty} a_\ell^2 = \frac{c\pi^2}{6}.$$

(3.6) <u>Theorem</u>. Let $\{Y(\theta):\theta \in E_k\}$ be as in (3.2) and let $\{a_\ell\}$ be a gain sequence. Then the following sequence of recursively defined random vectors

(3.7)   $\theta_{\ell+1} = \theta_\ell - a_\ell Y(\theta_\ell)$, $\theta_1$ arbitrarily chosen, converges in mean square and with probability one to $\theta$.

<u>Proof</u>. We refer the reader to the algorithm described in [1, pp 332-333] and the convergence proof given in the appendix to [1, pp 350-352]. Replacing their gain sequence $\{\rho_k\}$ with the gain sequence $\{a_\ell\}$, and replacing their norm $\|\cdot\|$ and inner product $<V,W> = V' W$ (where $V'$ denotes the transpose of the vector $V$ and $V' W$ denotes matrix multiplication) with the B-norm and B-inner product respectively, the theorem will follow once we verify that conditions (A1) - (A3) in [1, pp 332-333] are satisfied.

(A1) Since $E(Y(\theta)) = \theta - \overline{F}(\mathcal{g}(\theta))$, this result follows from (3.3) since B positive definite implies $\|\theta-\theta^\circ\|_B = 0$ if and only if $\theta = \theta^\circ$.

(A2) Follows from (3.3).

(A3) $E(\|Y(\theta)\|_B^2) = E(<\overline{Y}(\mathcal{g}(\theta)) - \mathcal{g}(\theta) + \theta, \overline{Y}(\mathcal{g}(\theta)) - \mathcal{g}(\theta) + \theta>_B)$

$\leq h(1 + \|\theta-\theta^\circ\|_B^2)$ for sufficiently large $h > 0$ since $S \subseteq D$ and by (2.3)

$\underset{\theta \in D}{\sup} E(\|\overline{Y}(\theta)\|^2) < \infty.$

## 4. Application to Maximum Likelihood Estimation

Let $D = \{(\alpha_i, \mu_i, \Sigma_i)\}_{i=1,\ldots,m}$ where each $\alpha_i > 0$ with $\sum_{i=1}^{m} \alpha_i = 1$, each $\mu_i \in R^n$ and each $\Sigma_i$ is a positive definite real symmetric $n \times n$ matrix. We consider $D$ to be realized as an open subset $D$ of $E_k$ where $k = \frac{m(n+1)(n+2)}{2} - 1$. For each $\theta = (\alpha_i, \mu_i, \Sigma_i)_{i=1,\ldots,m}$ let $x(\theta)$ be a random variable with values in $R^n$ and with distribution function

$$p(\theta, x) = \sum_{i=1}^{m} \alpha_i \, p_i(x) \qquad \text{for} \quad x \in R^n$$

where

$$p_i(\theta, x) = (2\pi)^{-n/2} |\Sigma_i|^{-1/2} \exp\{-\frac{1}{2}(x-\mu_i)^T \Sigma_i^{-1}(x-\mu_i)\}$$

for each $i = 1,\ldots,m$.

Fix $\theta^o \in D$ and let $\{x_k\}_{k=1,\ldots,N} \subseteq R^n$ be an independent sample of observations of $x(\theta^o)$. A maximum-likelihood estimate of $\theta^o$ based on $\{x_k\}$ is a choice of $\theta \in D$ which locally maximizes the log-likelihood function

$$L = \sum_{k=1}^{N} \log p(\theta, x_k)$$

In the appendix to [3], Peters and Walker prove there exists a sufficiently small neighborhood of $\theta^o$, such that with probability $\to 1$ as $N \to \infty$, there exists a unique maximum-likelihood estimate of $\theta^o$ in that neighborhood. Furthermore, with probability $\to 1$ as $N \to \infty$, this estimate $\to \theta^o$. This estimate will be called the consistent maximum-likelihood estimate (which we abbreviate by c.m.l.e.)

Equating $\nabla_\theta L = 0$ and performing algebraic manipulation of the resulting equations, yields the following necessary conditions for a m.l.e. $\theta$ of $\theta^o$.

$\theta = \mathcal{J}(\theta)$ where $\mathcal{J}:D \rightarrow D$ is defined by

$$\mathcal{J}((\alpha_i,\mu_i,\Sigma_i)_{i=1,\ldots,m}) = (\alpha_i',\mu_i',\Sigma_i')_{i=1,\ldots,m}$$

where, for each $i = 1,\ldots,m$

(4.1)
$$\alpha_i' = \frac{\alpha_i}{N} \sum_{k=1}^{N} x_k \frac{p_i(x_k)}{p(x_k)}$$

(4.2)
$$\mu_i' = \left\{\frac{1}{N} \sum_{k=1}^{N} x_k \frac{p_i(x_k)}{p(x_k)}\right\} \Big/ \left\{\frac{1}{N} \sum_{k=1}^{N} \frac{p_i(x_k)}{p(x_k)}\right\}$$

(4.3)
$$\Sigma_i' = \left\{\frac{1}{N} \sum_{k=1}^{N} (x_k - _i)(x_k - _i)^T \frac{p_i(x_k)}{p(x_k)}\right\} \Big/ \left\{\frac{1}{N} \sum_{k=1}^{N} \frac{p_i(x_k)}{p(x_k)}\right\}$$

where each $p_i$ and $p$ is evaluated with respect to the parameters $\theta = (\alpha_i,\mu_i,\Sigma_i)_{i=1,\ldots,m}$. $\mathcal{J}$ will be called the likelihood function.

In [3], Peters and Walker develop an iterative procedure which, starting with any initial estimate $\theta'$ which is sufficiently close to $\theta^o$, yields a sequence in $D$ converging to the c.m.l.e. of $\theta^o$ based on $\{x_k\}_{k=1,\ldots,N}$. Their technique consists in proving that, for $\varepsilon < \dfrac{4}{m(n+1)(n+2)}$, the function

(4.4)
$$\mathcal{J}_\varepsilon(\theta) = (1-\varepsilon)\theta + \varepsilon \, \mathcal{J}(\theta)$$

is locally contractive (at the c.m.l.e. of $\theta$) throughout a neighborhood of

$\theta^o$. Thus, for any $\theta'$ in this neighborhood of $\theta^o$, the sequence $\{\theta_\ell\}$ defined recursively by

(4.5)                $\theta_{\ell+1} = \mathcal{L}_\epsilon(\theta_\ell)$ , $\theta_1 = \theta'$

converges to the c.m.l.e. of $\theta$ .

In concluding, they discuss the computational advantages of this procedure as compared to classical numerical techniques such as Newton's method or the method of scoring. In particular, the procedure satisfies the following conditions.


(4.6)   At each stage of the iteration in (4.5), the constraints on the parameters in $\theta$ are satisfied.

(4.7)   The 'step size' $\epsilon$ depends only on $n$ and $m$ and not on $\theta^o$.

(4.8)   The procedure does not require the inversion, at each stage of the iteration, of a $k \times k$ matrix.

We will present a stochastic approximation analog of the iterative procedure defined by (4.5). In contrast to the classical stochastic method of scoring, our procedure satisfies the conditions (4.6) and (4.7). A step size $\epsilon$ will not appear explicitly in our algorithm.

Fix $\theta^o \in D$ and let $\{x_k\}_{k=1,\ldots,\infty}$ be an infinite sequence of independent samples of observations of $x(\theta^o)$. For any function $g$ from $R^n$ to any real vector space $V$, let

$$E(g) = \int_{R^n} g(x)p(\theta^o,x)dx$$

be the expectation, if it exists, of $g \circ x(\theta^{\sigma})$. ($\circ$ denotes composition of functions). Then, by the strong law of large numbers, with probability one as $N \to \infty$, the equations $(4.1) - (4.3)$ converge to

$$(4.9) \qquad \alpha_i' = \alpha_i E(\frac{p_i}{p})$$

$$(4.10) \qquad \mu_i' = E(x \frac{p_i}{p}) / E(\frac{p_i}{p})$$

$$(4.11) \qquad \Sigma_i' = E((x-\mu_i)(x-\mu_i)^T \frac{p_i}{p}) / E(\frac{p_i}{p}).$$

We denote the corresponding limiting value of the likelihood function by $\mathcal{I}$. Clearly $\mathcal{I}(\theta)$ is a continuously differentiable function of $\theta$ and $\mathcal{I}(\theta^{\sigma}) = \theta^{\circ}$; also by $(4.4)$, $\mathcal{I}_{\epsilon}(\theta) = (1-\epsilon)\theta + \epsilon \mathcal{I}(\theta)$ is locally contractive at $\theta^{\circ}$. By $(2.2)$, $\rho(\nabla \mathcal{I}_{\epsilon}(\theta^{\circ})) < 1$, implying that the eigenvalues of $\nabla \mathcal{I}(\theta^{\circ})$ have real parts strictly less than 1. Now define a function $d:D \to D$ by

$$(4.12) \qquad d((\alpha_i,\mu_i,\Sigma_i)_{i=1,\ldots,m}) = (\alpha_i,\alpha_i\mu_i,\alpha_i\Sigma_i)_{i=1,\ldots,m}.$$

Clearly $d$ is a differentiable function from $D \to D$ such that $d^{-1}$ exists and is differentiable. Define a function $\mathcal{I}':D \to D$ by $\mathcal{I}'(\theta) = d \circ \mathcal{I} \circ d^{-1}(\theta)$. By the chain rule for Frechet Derivatives,

$$\nabla \mathcal{I}'(d(\theta^{\sigma})) = [\nabla d(\theta^{\circ})][\nabla \mathcal{I}(\theta^{\circ})][\nabla d(\theta^{\circ})]^{-1}$$

hence the function

$$(4.13) \qquad \mathcal{J}'_\varepsilon(\theta) = (1-\varepsilon)\theta + \varepsilon \mathcal{J}'(\theta)$$

is locally contractive throughout some neighborhood $W$ of $d(\theta^o)$.

Define a family $\{\overline{Y}(\theta):\theta \in D\}$ of random variables with values in $E_k$ by

$$(4.14) \qquad \overline{Y}(\theta) = \varepsilon(\alpha_i - \alpha_i \frac{p'_i}{p'}, \; \mu_i - \alpha_i x(\theta^o) \frac{p'_i}{p'}, \; \Sigma_i - \alpha_i (x(\theta^o) - \frac{\mu_i}{\alpha_i})(x(\theta) - \frac{\mu_i}{\alpha_i})^T \frac{p'_i}{p})_{i=1,\ldots,m}$$

where $p'_i = p_i(d^{-1}(\theta), x(\theta^o))$ and $p' = p(d^{-1}(\theta), x(\theta^o))$.

Then $E(\overline{Y}(\theta)) = \theta - \mathcal{J}'(\theta)$. Therefore, the family $\{\overline{Y}(\theta):\theta \in D\}$ satisfies conditions (2.3) and (2.4).

Let $\{Y(\theta):\theta \in E_k\}$ be constructed from $\{\overline{Y}(\theta):\theta \in D\}$ as in (3.2) and let $\{a_\ell\}$ be any gain sequence. Then by (3.6), the sequence in (3.7) converges in mean square and with probability one to $d(\theta^o)$. Since $\{\varepsilon a_\ell\}$ is a gain sequence whenever $\{a_\ell\}$ is a gain sequence, $\varepsilon$ need not appear explicitly in the sequence in (3.7).